

**METHOD FOR PREDICTING FILE DOWNLOAD TIME FROM
MIRRORED DATA CENTERS IN A GLOBAL COMPUTER NETWORK**

5 This application is based on and claims priority from Provisional Application Serial
No. 60/208,013, filed May 26, 2000.

Related Applications

 This application is related to the following commonly-owned applications: "Global
Load Balancing Across Mirrored Data Centers," Serial No. xx/yyy,zzz, filed May 29,
10 2001; "Method For Generating A Network Map," Serial No. yy/xxx,zzz, filed May 29,
2001; and "Method For Extending A Network Map," Serial No. xy/xxx,yyy, filed May 29,
2001.

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2

BACKGROUND OF THE INVENTION

Technical Field

The present invention relates generally to high-performance, fault-tolerant content delivery and, in particular, to systems and methods for balancing loads from mirrored data centers within a global computer network.

Description of the Related Art

It is known to store Web-based content in mirrored data centers and to load-balance such content requests to data centers based on network traffic conditions. Existing global load balancing products use several different approaches for building a map of Internet traffic conditions. One approach uses border gateway protocol (BGP) data. BGP-based routing, however, can be sub-optimal because the BGP data can be very coarse. Other approaches attempt to compute an optimal mapping in real-time and then cache the mapping information. This technique can lead to poor turnaround time during an initial “hit” and potentially stale mappings on successive requests. In addition, the quality of the measurement to the endpoint tends to be noisy. Because of the deficiencies of these mapping techniques, the resulting load balancing is less than effective.

Current load balancing devices are typically incapable of computing an optimal map for an entire computer network such as the entire Internet. Presently, the Internet has hundreds of millions of hosts and routers. Estimating the connectivity time of the entire Internet to a set of mirrored data centers, such as by evaluating the network path between a server and each and every host or router, would be incredibly time-consuming and would consume far too much bandwidth. Such techniques, of course, are impractical when real-time routing decisions are required.

Further, such measurements tend to be noisy and inaccurate, and they can annoy system administrators whose firewalls are contacted. Local name servers behind firewalls would not be reached and slow connectivity over the “last mile” (e.g., due to dial-up connections and the like) tend to confuse the connectivity picture. Consequently, there remains no efficient technique in the prior art for generating an optimal network connectivity map that can be used for providing intelligent traffic redirection in conjunction with load balancing across mirrored data centers located around the globe.

BRIEF SUMMARY OF THE INVENTION

The invention is an intelligent traffic redirection system that does global load balancing. It can be used in any situation where an end-user requires access to a replicated resource. The method directs end-users to the appropriate replica so that the route to the
5 replica is good from a network standpoint and the replica is not overloaded. The technique preferably uses a Domain Name Service (DNS) to provide IP addresses for the appropriate replica.

In a preferred embodiment, the method relies on a network map that is generated continuously for the user-base of the entire Internet. The problems inherent in the prior art
10 are overcome by vastly reducing the dimensionality of the problem of estimating the relative connectivity to a set of mirrored data centers. A "data center" is typically located at a telecommunications facility that leases space and sells connectivity to the Internet. Multiple content providers may host their web sites at a given data center. Instead of probing each local name server (or other host) that is connectable to the mirrored data
15 centers, the network map identifies connectivity with respect to a much smaller set of proxy points, called "core" (or "common") points. A core point then becomes representative of a set of local name servers (or other hosts) that, from a data center's perspective, share the point. Each set of mirrored data centers has an associated map that identifies a set of core points.

Once core points are identified, a systematic methodology is used to estimate
20 predicted actual download times to a given core point from each of the mirrored data centers. According to the invention, ICMP (or so-called "ping" packets) are used to measure roundtrip time (RTT) and latency between a data center and a core point. Thus, for example, a core point may be pinged periodically (e.g., every 30 seconds) and the
25 associated latency and packet loss data collected. Using such data, an average latency is calculated, preferably using an exponentially time-weighted average of all previous measurements and the new measurement. A similar function is used to calculate average packet loss. Using the results, a score is generated for each path between one of the data centers and the core point. The score may be generated by modifying an average latency,
30 e.g., with a given penalty factor, that weights the average latency in a unique way to

provide a download prediction. Whichever data center has the best score (representing the best-performing network connectivity for that time slice) is then associated with the core point.

The score is a good approximation of the amount of time that may be required to download an average size test file from one of the mirrored sites to a core point.

According to another technical advantage of the present invention, a method of predicting a file download time begins by periodically initiating a test probe from a server to a given point in a network. Latency and packet loss data generated from the test probes is then collected. This data is then used to compute an exponentially time-weighted average of latency and a time-weighted average of loss. The score is then generated as a function of the time-weighted average of latency modified by a penalty factor that is a function of the time-weighted average of loss.

The generated network map is then used to effect traffic redirection and load balancing. In particular, when a user's local name server makes a request for the content provider's web site (located within a set of mirrored data centers), the method preferably uses the network map to return to the local name server a list of web server IP addresses at the optimal data center. If ping data is not available for the user's local name server (of it the IP block has not been extended through unification), BGP or geo-routing can be used to make a default routing decision. Content provider-specified load balancing preferences may also be readily enforced across the data centers and/or within a particular data center.

The foregoing has outlined some of the more pertinent objects and features of the present invention. These objects should be construed to be merely illustrative of some of the more prominent features and applications of the invention.

BRIEF DESCRIPTION OF THE DRAWINGS

For a more complete understanding of the present invention and the advantages thereof, reference should be made to the following Detailed Description taken in connection with the accompanying drawings, in which:

5 Figure 1 is an illustration of a mirrored Web site that is managed by a global traffic manager according to the present invention;

 Figure 2 is a high level illustration of the components of the GTM service;

 Figure 3 is a simplified illustration of a core point discovery process of the invention;

10 Figure 4 is a simplified illustration of how an end user request is processed by the global traffic redirection system of the present invention for a mirrored web site that has been integrated into the managed service;

 Figure 5 is a flowchart describing how a map is generated by the GTM system;

15 Figure 6 is a simplified block diagram of one implementation of the global traffic management system of the invention; and

 Figure 7 is a representative traceroute generated during the core point discovery process.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

By way of brief background, it is known in the art for a Web content provider to distribute or “mirror” its Web site to ensure that the site is always available and providing acceptable performance for a global customer base. Once a Web site is distributed, global traffic management (GTM) solutions typically are used to direct users to the various mirror sites. GTM solutions use a variety of methods to determine which is the “best” mirrored site in which to direct a user. Because Internet conditions are constantly changing, however, the “best” site for a particular user also varies with these conditions. The present invention is a GTM solution that maximizes availability and performance of a mirrored delivery site.

In a preferred embodiment now described, the global traffic management solution is a managed service provided by a service provider, such as a content delivery network (CDN) service provider (CDNSP). As is well-known, a CDN is a network of geographically distributed content delivery nodes that are arranged for efficient delivery of digital content (e.g., Web content, streaming media and applications) on behalf of third party content providers. A request from a requesting end user for given content is directed to a “best” replica, where “best” usually means that the item is served to the client quickly compared to the time it would take to fetch it from the content provider origin server. Typically, a CDN is implemented as a combination of a content delivery infrastructure, a request-routing mechanism, and a distribution infrastructure. The content delivery infrastructure usually comprises a set of “surrogate” origin servers that are located at strategic locations (e.g., Internet Points of Presence, network access points, and the like) for delivering copies of content to requesting end users. The request-routing mechanism allocates servers in the content delivery infrastructure to requesting clients in a way that, for web content delivery, minimizes a given client’s response time and, for streaming media delivery, provides for the highest quality. The distribution infrastructure consists of on-demand or push-based mechanisms that move content from the origin server to the surrogates. An effective CDN serves frequently-accessed content from a surrogate that is optimal for a given requesting client. In a typical CDN, a single service provider operates the request-routers, the surrogates, and the content distributors. In addition, that service

provider establishes business relationships with content publishers and acts on behalf of their origin server sites to provide a distributed delivery system. A well-known commercial CDN that provides web content and media streaming is provided by Akamai Technologies, Inc. of Cambridge, Massachusetts.

5 Thus, in one embodiment, the present invention implements a managed service for global load balancing of a content provider's mirrored Web sites. Figure 1 illustrates the basic implementation environment. In this example, the global traffic management service 100 provides global traffic management for a content provider running a pair of mirror Web sites 102 and 104 (identified by the same domain, e.g., www.akamai.com). The
10 GTM service 100 provides improved responsiveness for end users 106 and 108 accessing the Web site by directing them to the best performing mirrored site. Figure 2 illustrates the high level technical architecture of the GTM service which, as noted above, is implemented by a CDNsp or other entity (the "managed service provider") as a managed service on behalf of content providers running mirrored Web sites. Of course, one of
15 ordinary skill will appreciate that the inventive functionality may also be implemented in whole or in part as a product-based solution.

For illustrative purposes only, and with reference to Figure 2, a preferred GTM service 200 comprises a number of components: a set of network agents 202, a set of web server agents 204, a set of map generation servers 206, and a set of name servers 208.
20 Each such component typically is a server, such as a Pentium-based box running the Linux operating system and having application software for carrying out the functions described below, or one or more processes executing on such a machine. As will be described, data is collected by the network agents and the web server agents and delivered to the map generation servers. The map generation servers analyze the data, and at least one map
25 server produces a map that assigns name server IP address/blocks to regions. At least one map is then uploaded to the name servers. When an end user requests access to a mirrored site domain being managed by the service, one of the name servers hands back an IP delegation that represents a "best" data center to which the user should connect.

The particular placement of the components as illustrated in the drawing is
30 representative, and there is no requirement that any particular entity own or control a

particular machine. In this embodiment, a content provider has network agents located in or near their network segment within each respective data center that hosts the mirrored Web site. Thus, for example, a pair of network agents 202a and 202b are dedicated to the content provider in data center 203a, and a pair of network agents 202c and 202d are
5 dedicated to the content provider in data center 203b, although this is not required. These network agents preferably share the same network connection as the content provider's web servers. Thus, e.g., network agents 202a and 202b in data center 203a share network connections with the content provider's web servers 207a-c. Where the managed service provider operates a CDN, the set of network agents may be deployed in data centers in
10 which the CDN is deployed. Of course, multiple content providers may host their web sites at a given data center and share network agents. Thus, a given network agent may collect data once for a first content provider at a given location and then share the data across all other content providers co-located in the same data center. A data center typically is located at a telecommunications facility (e.g., Exodus, Frontier Global, UUUNet, and the like) that leases space and sells connectivity to the Internet.

A network agent has two (2) primary functions: running "core point" discovery (CPD) to determine a set of "core" points, and monitoring network performance to each core point. As will be seen, the inventive system continuously pre-computes optimal maps, preferably for the user base of the entire Internet. It is able to do this effectively
20 because the system reduces the scale of the problem by aggregating parts of the Internet and representing them with "core" points. A core point typically is representative of a set of local name servers (or other hosts) that, from the perspective of a given network location (e.g., a data center), share the point. Typically, a core point is a router on the Internet, although this is not a requirement. The information collected from the core point
25 discovery process is fed to the map generation servers on a relatively frequent basis, e.g., one every thirty (30) seconds, to identify down routes, congestion, route changes, and other network traffic conditions that may impair or effect connectivity to a data center at which a particular mirrored site is hosted.

According to a preferred embodiment of the invention, a core (or "common") point
30 is discovered as follows. An incremental trace route is executed from each of the set of

mirrored data centers to a local name server that may be used by client to resolve a request for a replica stored at the data centers. An intersection of the trace routes at a common routing point is then identified. Thus, for example, the common routing point may be the first common point for the trace routes when viewed from the perspective of the data centers (or the last common point for the trace routes when viewed from the perspective of the local name server). The common routing point is then identified as the core point for the local name server. A core point is identified for other local name servers (or other hosts) in the same manner. Thus, a given set of mirrored data centers may have associated therewith a set of core points that are then useful in estimating the relative connectivity to the set of data centers, as is described below.

Figure 3 is a simplified diagram of the core point discovery process, in accordance with one embodiment of the invention. For purposes of example only, in Figure 3, the data center 300 corresponds to a data center located on the West Coast and the data center 302 corresponds to a data center located on the East Coast. Data center locations, of course, are merely representative. Each data center can host a mirror site for a given content provider. According to the invention, a core point 305 is discovered as follows. An incremental trace route is executed from each of a set of mirrored data centers 300, 302 to local name servers 304, 306, 308 that may be used by a client machine 310. For example, in Figure 3, the network agent (not shown) has executed a first set of traceroutes, between the data center 300 and the local name servers 304, 306 and 308, and a second set of traceroutes between the data center 302 and the local name servers 304, 306 and 308. The network path between the respective data center and the local name server(s) contain router hops, as is well known. To locate a core point, the network agent identifies a location at or close to the intersection of the trace routes at a common routing point, which is shown in Figure 3 as a core point 305. For example, the common routing point may be the first common point for the trace routes when viewed from the perspective of the data centers 300 and 302 (or the last common point for the traceroutes when viewed from the perspective of the local name server 304). The common routing point is then identified as the core point 305 for the local name server. Figure 8 illustrates a representative core point discovery process trace.

For example, if two or more different paths are traced and the same route (or routes) appears on at least a portion of all of the paths, the common routing point can lie somewhere along that common portion of the route. As noted above, generally the core point is the first common point for the trace routes when viewed from the perspective of the data centers, which is the same as the last common point for the trace routes when viewed from the perspective of the local name server.

The core point 305 need not be situated at the “exact” intersection of the trace routes. It can, for example, be located near or substantially near the intersection. It can also be located adjacent to the intersection, or it can be located at any nearby point such that measurements made to the point are representative of the measurements made at the intersection.

The network agent identifies other core points for other local name servers (or other hosts) in the same manner. Thus, a given set of mirrored data centers may have associated therewith a set having one or more core points that are then useful in estimating the relative connectivity to the set of data centers, as is described below. If network paths on the Internet are changing frequently, a network agent preferably runs core point discovery with some frequency.

As noted above, a network agent also performs the function of periodically checking the core points assigned to one or more local name servers that already have been mapped. This process is now described.

Network agents preferably make measurements to core points using Internet Control Messaging Protocol (ICMP) (or so-called “ping” packets) to evaluate such information as round trip times (RTTs), packet loss, and number of router hops. Thus, using the example in Figure 3, a given network agent periodically “pings” a core point (e.g., every 30 seconds) and collects the associated latency and packet loss. Using such data, the network agent calculates an average latency. In one embodiment, the network agent calculates average latency using an exponentially time-weighted average of all previous measurements and the new measurement. The network agent uses a similar function to calculate average packet loss. This calculation is described in more detail below. Using the results, the network agent generates a “score” for each path between one

of the data centers and the core point. The score is generated, for example, by modifying an average latency with a given penalty factor that weights the average latency in a unique way to provide a download prediction. Whichever data center has the best score (representing the best-performing network connectivity for that time slice) is then associated with the core point.

Referring back to Figure 2, the web server agents 204 do test downloads to either all the web server IP addresses or to the local load balancing devices to test for availability or “aliveness” of the mirrored sites (i.e., per data center mirror or web server). Typically, a web server agent tests an object, e.g., a twenty (20) byte file available on the web server via an HTTP GET request, and check for errors and download times. In a representative embodiment, the measurements are taken periodically, e.g., every ten (10) seconds, although preferably a customer can change the timeout. An IP address is declared “dead” if more than a given percentage of the web server agents are unable to download the test object within the timeout threshold. This allows customers to set a threshold on response times so that the system can direct traffic away from data centers where performance suffers. The web server agents are preferably dispersed in co-location facilities, which are dispersed geographically and on a network basis. Moreover, one skilled in the art will recognize that the described functions of the web server agent could be performed by another component, such as the network agent, the map generation server, or some other server. Moreover, neither the web server agent nor its functions (such as testing the aliveness of a data center) are necessary for certain embodiments of the invention.

The map generation servers 206 receive data from the network agents and the web server agents and use this data to generate maps, which describe the mirrored site that is optimal for each IP address block. In a preferred embodiment, a map is achieved by evaluating web server agent data, a time-weighted average of latency and packet loss, and BGP and geo information. Preferably, there are two (2) map generation server processes for each customer, and maps are generated periodically, e.g., every 3-5 minutes. Although not a limitation, preferably the map generation servers associate IP blocks with Internet “regions” such that a given map associates an IP block with a region number. Another data file is then used to associate region number to physical IP address. In a representative

embodiment, maps (which associate IP block to region #) are generated every few minutes and then uploaded to the name servers.

The name servers 208 hand out to the requesting end user the IP address(es) of the optimal data center. Typically, the name server response have a time to live (TTL) of about five (5) minutes, although this value may be customer-configurable. In a representative embodiment, the name servers are the same name servers used by the CDNSP to facilitate routing of end user requests to CDN content servers.

Figure 4 illustrates how a customer web site is integrated into the traffic redirection system of the present invention. In a representative embodiment, it is assumed that the customer has a distributed web site of at least two (2) or more mirrored sites. The inventive system load balances multiple subdomains/properties provided they are in the same data centers. Integration simply requires that the customer set its authoritative name server 400 to return a CNAME to the GTM name servers 408, which, thereafter, are used to resolve DNS queries to the mirrored customer site. Recursion is also disabled at the customer's authoritative name server. In operation, an end user 402 makes a request to the mirrored site using a conventional web browser or the like. The end user's local name server 404 issues a request to the authoritative name server 400 (or to a root server if needed, which returns data identifying the authoritative name server). The authoritative name server then returns the name of a name server 408 in the managed service. The local name server then queries the name server 408 for an IP address. In response, the name server 408 responds with a set containing one or more IP addresses that are "optimal" for that given local name server and, thus, for the requesting end user. As described above, the optimal set of IP addresses is generated based on network maps created by testing the performance of representative core points on the network. The local name server selects an IP address from the "optimal" IP address list and returns this IP address to the requesting end user client browser. The browser then connects to that IP address to retrieve the desired content, e.g., the home page of the requested site.

Figure 5 is a high level flowchart illustrating how data is processed in order to create a map. Periodically (e.g., every thirty (30) seconds), the network agents ping each core point from each data center. This is step 500. At each network agent, a time-

weighted average of latency, and a time-weighted average of loss, is computed. This is step 502. As will be described, the weights decay exponentially in time with a time constant that is configurable. At step 504, the data is further processed to produce a score for each data center per core point. At step 506, each core point is then associated with the name servers for which the core point was a proxy. At step 508, a map generation process goes through all of the data and decides a set of candidate data centers for each name server. At this time, any data centers that the web server agents determine are not “alive” are discarded. At step 510, the map generation process extends its decisions with respect to name servers to decisions with respect to IP block. A unifying algorithm is used to provide this functionality. This algorithm operates generally as follows. If all name servers in a BGP-geo block have agreeing ping decisions, then the decision of what data center is “optimal” is applied to the whole block. Conversely, if there is a disagreement, the block is broken up into the largest possible sub-blocks so that, in each sub-block, all the name servers agree. For any block that has no name servers, the BGP-geo candidates may be used.

Referring now back to Figure 5, at step 512, the map is produced with the candidate for each block. If there are multiple candidates, the assignments are made to get as close to the load balancing targets as possible. The load balancing targets are defined, usually by the content provider, and these targets may be percentages (adding up to 100%) that breakdown the desired traffic amount by data center. This completes the map generation process.

As described above, step 502 involves generating a time-weighted average of latency and a time-weighted average of loss. More generally, this aspect of the invention provides a systematic methodology for predicting actual download times for various flow control protocols, e.g., TCP. As is known, TCP is the most commonly used flow control protocol on the Internet today. Other protocols are built on top of UDP. Neither TCP nor UDP packets can be used to monitor the state of routes on the Internet, however. According to the present invention, ICMP packets are injected into the network (e.g., by the network agents), at preferred points in time, and then routed to a suitably chosen intermediate core point. The system then looks at the behavior of the Internet induced by

the ICMP probes by computing latency and packet loss. Latency is a measure of the round trip time (RTT) between the server and the core point. From maintaining a time series of loss and latency, the system is able to predict effectively the amount of time it would take a client (that uses a name server associated with the core point) to initiate and complete a download from the server. The quality of this prediction is important for effective mapping because when a client makes a web request and there are multiple web servers from which to potentially server, it is important to be able to predict correctly which web server has the best connectivity. This is a difficult problem in general because the Internet is highly bursty and exhibits highly variable traffic conditions.

The following example illustrates how the time-weighted averages are computed in accordance with one embodiment of the invention. Assume for purposes of example only that a content provider (Figure 3) has mirror sites located at two data centers 300 (West Coast) and 302 (East Coast). The network agent “pings” the core point 305 from each data center. The network agent stores the latency and the packet loss for each measurement made. It should be understood that latency and loss parameters are merely representative of the types of signal transmission parameters that the network agent can track. Other parameters that could be measured include any parameter helpful in determining the speed, quality and/or efficiency of a network path, such as parameters indicative of outages on paths, loss in signal strength, error-control data, route changes, and the like.

For each “ping” to/from each data center to the core point, the respective network agent logs the data. Table 1 illustrates an example of the type of data that the network agent gathers over the course of measurements made every 30 seconds between the data centers and the core point. Table 1 is a table of latency measurements (data is in seconds (s)) and shows the current measurement (t=0) followed by measurements made previously.

Table 1

Parameter	Data Center	Current	t-30s	t-60s	t-180s	t-240s	t-300s	Avg (s)
Latency	(West)	8.0	7.5	7.7	8.2	7.6	7.7	7.78
	(East)	3.0	3.5	3.2	3.8	3.6	3.4	3.42
Loss	(West)	0	0	1	0	1		N/A
	(East)	0	0	0	0	0		N/A

As Table 1 shows, based on latency, in this example the East Coast data center appears to have a smaller average latency to the core point than the West Coast data center.

- 5 A time-weighted average of latency, and a time-weighted average of loss, is then computed. The weights decay exponentially in time with a time constant that is configurable (e.g., a time constant of 300 seconds). For a sequence of measurements made (t_i, x_i) , where t_i is the time of the i^{th} measurement and x_i is the value measured (e.g., x_i can be the latency measurement lat_i or the loss measurement $loss_i$), the time weighted average of latency is computed as:

$$AverageLatency = \sum_{i=0}^{\infty} lat \times e^{-t_i / C}$$

Assuming that the time constant $C = 300$ seconds, and using the data of Table 1, the average latency time series is computed as:

$$AverageLatency = \sum_{i=0}^{\infty} lat \times e^{-t_i / 300}$$

- 15 Using the data, the average latency for the data center 300 is computed as:

$$AverageLatency = \sum_{i=0}^{\infty} (8.0e^{0/300} + 7.5e^{-30/300} + 7.7e^{-60/300} + 8.3e^{-180/300} + 7.6e^{-240/300} + 7.7e^{-300/300})$$

$$AverageLatency = \sum_{i=0}^{\infty} (8.0(1) + 7.5(.9048) + 7.7(.8187) + 8.3(.5488) + 7.6(.4493) + 7.7(0.3678))$$

$$AverageLatency = \sum_{i=0}^{\infty} (8.0 + 6.78 + 6.31 + 4.55 + 3.41 + 2.83)$$

$$AverageLatency = \sum_{i=0}^{\infty} (31.88)$$

To compute the exponentially time weighted average, the network agent sums each weighted latency measurement (e.g., 31.88) and divides this sum by the sum of the weight factors (i.e., $e^{-30/300} + e^{-60/300} \dots$ etc.). Thus, the exponentially time weighed average latency for the data center 300 is computed as:

5 Exponentially time weighted average = $31.88/4.0894$
 Exponentially time weighted average = 7.795

As these computations show, the exponentially time-weighted average is 7.79, which differs from the computed average of 7.78. Although this difference does not appear significant in this example, it can be more significant for measurements averaged
10 out over long periods of time, because more recent measurements will be given more weight than older measurements. The network agent determines dynamically whether core points that were once considered optimal are still so, whether core points that had been performing well (for a given time period) are now degraded, and the like. The exponentially time weighted averaging helps also to smooth out aberrations over time in
15 measured data and helps to indicate trends.

Using the same information, the time weighted average latency for the East Coast data center 302 are computed in a similar manner. In addition, although not illustrated here, the network agent computes a time-weighted average of loss in the same way.

As described above, time-weighted averages are then processed to produce a score
20 for each data center per core point. A preferred scoring function is as follows:

Score function = average latency + {[max (100, average latency)]*(penalty factor)},

where the score is a value in milliseconds. Each of the values has a millisecond unit, except for the penalty factor, which is unit-less. The value “100” is a floor or base-level
25 value, which represents the usual round trip time required for a packet to travel between the East Coast and the West Coast. The floor is variable. The term “max” refers to selecting either the value “100” or the average latency, whichever is greater. That value is then multiplied by a penalty factor, with the result then being added to the average latency to generate the score. The penalty factor preferably is a function of the time weighted
30 average loss. Thus, in one illustrative embodiment, the penalty factor is some multiple of

the time weighted average loss. The multiplier may be varied, e.g., as a function of percentage of loss, with the penalty factor contribution being higher for greater packet loss. In a given embodiment, the scoring function may have the following variants:

For losses less than a given percentage (e.g., 10%), the scoring function is
5 computed as:

$$\text{Score} = \text{average latency} + \{[\max(100, \text{average latency})] * (10 * \text{average loss})\}$$

For losses greater than the given percentage, the score is computed as:

$$\text{Score} = \text{average latency} + \{[\max(100, \text{average latency})] * ((110 * \text{average loss}) - 10)\}$$

10 The penalty factors in the scoring function variants are merely representative.

According to the invention, it has been found that a scoring function such as described above that is based on time-weighted average latency weighted by a time-weighted average loss penalty factor affords a good approximation or “proxy” of the download time for an average size (e.g., 10Kbyte) file from the data center to an average
15 end user. Of course, the file download time would be expected to vary as the file size is varied, but it has been found that the scoring function described above still tends to capture which data center of the mirrored set provides better performance. In other words, the absolute value of any given score is not as important as the data center-specific (e.g., East Coast vs. West Coast) values.

20 When the scores are provided to the map generation process, the network agent associates the core point with the local name server(s) for which the core point serves as a “proxy.”

The following describes a specific implementation of the global traffic redirection system as a managed service offering on behalf of content providers running mirrored web
25 sites. Figure 6 illustrates the overall system architecture 600. As noted above, these processes typically run across multiple servers in the system. There are three logical grouping of these processes. First, the PingServer 602, PingProcessor 604, and TestPingServer 606 are running on the network agents located in the content provider’s data centers. Second, the MapMaker 608, MapTester 610, and DBPusher 612 are running
30 on another set of servers. However, these may also be run on the network agent machines

if there is a shortage of servers in the network in which the global traffic management system operates. Another set of processes, called MapNote Web 614 and MapNoteDNS 616, run together on a relatively static set of machines for all customers of the system. Processes 602, 604, 608, 610, 612, 614 and 616 typically run continuously. An alert processor (not shown) detects if one or more machines on the network are non-functional and sends one or more corresponding alerts. An archive process (not shown) is used to automatically log files and other system files periodically. A file moving process (not shown) is used move data files. Each server may also run a generic process monitor (not shown), which reports data to a service provider query system.

As has been described, the global traffic management system 600 collects several pieces of data that results in a map being uploaded to the GTM name servers 615. At the beginning, Core Point Discovery (CPD) produces a list of IP addresses in a file (signifying the core points). This file is sent to each PingServer 602. Preferably, there is a PingServer process 602 running on each of the network agents that are deployed in a content provider's data center (not shown). In this embodiment, there is a pair of machines in each data center, only one PingServer process is primary. The other one is running but only takes over if the primary goes down. Each PingServer process 602 pings each of the core points approximately every 30 seconds.

Next, the ping results are sent to the PingProcessors 604. PingProcessors 604 preferably run on the same machines as the MapMakers 608, although this is not a requirement. The PingProcessors 604 process the ping results and drop the data files off for the MapMaker 608. MapMakers 608 also require data from the MapNoteWeb agents 614. The MapNoteWeb agents 614 are the web server agents that do test downloads from the content provider's web servers. These tests are used to determine aliveness of the web servers in the data centers as has been described.

The MapMaker 608 looks at the ping data as well as the MapNote Web data and creates a top-level map for the top-level name servers. The map is then sent to the MapTester 610 (which is usually running on the same machine). The MapTester 610 uses test ping data from the TestPingServer 606 to check a given number of (e.g., a few hundred) IP addresses in the map. This is done to make sure the map is correct, however,

this processing is optional. Finally, if the map passes the test, it is queued for uploading to the name servers 615.

DBPusher 612 is one other process that preferably runs on the same machines as the MapMaker process 608. This process is solely responsible for pushing out a DB file to the top-level name servers 615. This DB file completes the lookup data for the top-level name server 615. That is, the map from the MapMaker 608 contains a mapping of IP block to a virtual region number. The DB file is the data file that actually has a mapping of the region number to physical IP addresses. DBPusher 612 monitors the MapNote Web data and, in case of a failure, pushes an updated DB file to the name servers.

PingServer 602 is responsible for measuring RTT and packet loss to the list of core points. The list of core points determined as follows. Preferably, there is a PingServer process running for each content provider at each data center in which a content provider is co-located. Thus, in one embodiment, the service provider deploys servers in all of a content provider's data centers. In another embodiment (not shown), ping data is shared for all customers who co-locate at a particular data center, and the GTM service provider may simply pre-deploy servers at "popular" hosting facilities to save time in integrating new customers to use the system.

The PingServer process preferably is run on each of the network agents in a data center. A leader election process (not shown) may be used to allow for the non-leader to take over if the primary fails within that same data center. PingServer includes a process that is used to ping a list of IP addresses, which the PingServer receives from a system source. Also, before the list is pinged, any IP addresses that are on a restricted list are filtered out. In particular, the primary inputs to the PingServer process are as follows:

- Restricted tree - a list of IP addresses that are not pinged.
- Routers file - the list of IP addresses that were discovered using Core Point Discovery.

The outputs of PingServer are as follows:

- Ping results - raw results of pinging IP addresses.
- Routers file - list of all IP addresses that PingServer used

PingProcessor is responsible for taking the raw measurement data from PingServer and computing the time-weighted averages. The time-weighted average is computed both for RTT and packet loss measurements for the core points. The time-weighted average is computed as described above. The primary inputs to the PingProcessor process are as

5 follows:

- Ping results from PingServer
- Routers file from PingServer

The outputs of PingProcessor are as follows:

- 10
- Nameserver list
 - Processed ping data

The MapMaker creates the map for the top-level name servers. MapMaker takes the processed ping data from PingProcessor and the aliveness data from MapNoteWeb and constructs a map. This map contains a relationship between certain IP blocks and a region number. The inputs to MapMaker are:

- 15
- Nameserver list from PingProcessor
 - Ping scores from PingProcessor
 - BGP-Geo tree information

The outputs of MapMaker may include, for example:

- 20
- Debug map
 - Map states file
 - Map
 - Ping data

MapTester is the last process before a map is uploaded to the top-level name servers. MapTester receives a candidate map from MapMaker. It then looks-up the mapping of a test IP addresses (that are pinged using TestPingServer, which is discussed more fully below). If the number of differences is below some threshold, then the map is deemed acceptable.

The map is then uploaded to the top-level name servers. The inputs to the MapTester process are:

- Debug map
- Test ping data
- Stats file
- Map
- Ping data

5

The output of the MapTester process is:

- Map

TestPingServer collects RTT and packet loss information for a small subset of IP addresses. This data is collected to ensure that the map produced by MapMaker is valid.

10

MapTester, assuming the map is good, will then load the maps to the top-level name servers. The inputs to the TestPingServer process are:

- Restricted tree
- List of IP addresses to test

The output of the TestPingServer process is:

- Test ping data

15

As noted above, because the MapMaker map only provides a mapping between IP block and a region number, a separate process preferably is used to provide the mapping between region number and the actual IP addresses of the webserver(s). DBPusher is responsible for processing the MapNoteWeb data and creating a DB file that is uploaded to the top-level name servers. Then, the top level name server will, after it has determined the region number for a give IP in the map, look in the corresponding DB file to find the right set of IP addresses to hand out. The input to DBPusher is:

- MapNote Web data

The output to DBPusher is

25

- DB file for name servers - this file is pushed to the name server directly by DBPusher

MapNote Web is run on a select number of servers for all customers using the traffic management system. MapNoteWeb uses a list of target URLs (which, for example, could be stored in its configuration files) on which it performs test downloads. Preferably,

these tests are simple HTTP GET requests and only time and errors are important. This data can be interpreted to determine whether or not a data center or web server is alive or dead. The download time is stored historically using a time-weighted average. Both the MapMaker and DBPusher use this data. The input to MapNoteWeb is:

- Targets to measure against (stored in configuration file)

The output to MapNoteWeb is:

- Download test results

MapNoteDNS is related to MapNoteWeb, except that instead of monitoring web servers it monitors name servers. Its basic function is to do a measurement at a name server for certain names and report the time. If no answer comes back, the process will attempt to ping to determine whether it is the name server or the network that caused the failure. The inputs to MapNoteDNS are:

- Name servers to test
- What domains to test for

The output of MapNoteDNS is:

- DNS query results

Although not described in detail, various tests (that are not relevant to the present invention) may be executed to determine whether or not each of the above-described processes is running correctly.

The intelligent traffic redirection system has numerous advantages. The system continuously pre-computes optimal maps for the user-based of the entire Internet (or, if desired, a given sub-portion thereof). It is able to do this effectively because the system reduces the scale of the problem by aggregating parts of the Internet and representing them with core points. The system is also able to make different kinds of measurements depending upon the service being replicated. It combines these measurements for the core points into decisions which it then extends to the entire Internet using unification over a fallback partition of the IP address space using, e.g., BGP and geo information. The system also is unique in its ability to balance load for cost minimization.

The system is able to pre-compute an optimal mapping for the entire Internet at all points in time. In addition to being extremely fast in its ability to react to bad network

conditions, it is also extremely fine-grained in its response. The system is able to detect bad server conditions quickly and is capable of interfacing with a multitude of local load balancers. By computing core points, the system makes superior measurements that mitigate the problem of intruding on firewalls and other detection mechanisms. Moreover, unlike the prior art, it can load balance load so as to minimize bandwidth costs.

Predicting download times using ICMP probes and time-series techniques also provides numerous advantages. The technique does not have any restriction on the range of file sizes and download types, and it makes intelligent use of ICMP probes of different sizes to effectively estimate packet loss. The technique requires very little state for keeping track of the time-series and is able to quickly compute a new estimate using an exponentially time-weighted average of all previous measurements and the new measurement. Rather than attempting to probabilistically model TCP flows, the inventive technique provides a general method for extracting a good predictor of download times based on ICMP probes.

In the preferred embodiment, the intelligent traffic redirection system is used to direct traffic to a mirrored Web site. Generalizing, the inventive system and managed service can be used in any situation where an end-user requires access to a replicated resource. As described above, the system directs end-users to the appropriate replica so that their route to the replica is good from a network standpoint and the replica is not overloaded. An "end user" may be generalized to any respective client system, machine, program or process. Thus, other uses of the system may include, without limitation, to direct caches to storage servers, to direct streaming servers to signal acquisition points, to direct logging processes to log archiving servers, to direct mail processes to mail servers, and the like.

Having thus described our invention, the following sets forth what we now claim.